



(19)

(11) Publication number:

0

Generated Document.

PATENT ABSTRACTS OF JAPAN(21) Application number: **05050398**(51) Intl. Cl.: **G06F 3/06 G06F 3/06**(22) Application date: **11.03.93**

(30) Priority:	(71) Applicant: HITACHI LTD
(43) Date of application publication: 22.09.94	(72) Inventor: FUJII TETSUHIKO YAMAMOTO AKIRA SATO TAKAO YOSHIDA MINORU
(84) Designated contracting states:	(74) Representative:

**(54) DISK ARRAY
CONTROL METHOD**

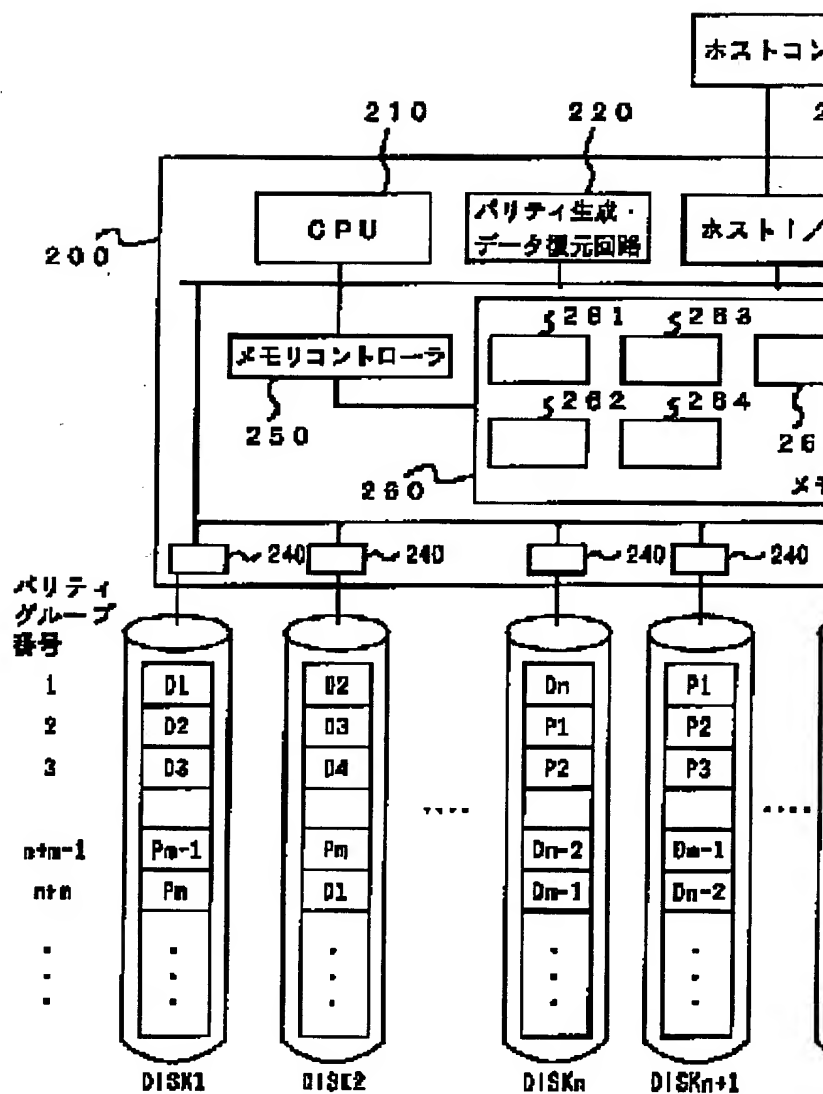
(57) Abstract:

PURPOSE: To prevent the increase of the frequency in access by using data and redundant data, which are stored in disk devices other than a disk device where a fault is detected, to restore data of this disk device and recording this data in the disk.

CONSTITUTION: If a parity group of a processing object is not restored yet, a classification (c) of data/redundant data stored in a faulty disk is obtained in this parity group. When the classification (c) indicates data, numbers of disks where data and redundant data belonging to the same parity group are stored are obtained from a parity constitution management table 264 for the purpose of restoring data, and they are read out from pertinent disks to restore data in the faulty disk. Restored data or generated redundant data is written over redundant data

selected for overwrite, and restoration state management information corresponding to the entry in the parity group of the processing object in a restoration state management table 263 is changed from the unrestored state to the already restored state to terminate the processing.

COPYRIGHT: (C)1994,JPO&Japio



(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平6-266508

(43) 公開日 平成6年(1994)9月22日

(51) Int.Cl.⁵

G 0 6 F 3/06

識別記号 庁内整理番号

3 0 4 B 7165-5B

3 0 5 F 7165-5B

F I

技術表示箇所

審査請求 未請求 請求項の数11 OL (全 27 頁)

(21) 出願番号 特願平5-50398

(22) 出願日 平成5年(1993)3月11日

(71) 出願人 000005108

株式会社日立製作所

東京都千代田区神田駿河台四丁目6番地

(72) 発明者 藤井 哲彦

神奈川県川崎市麻生区王禅寺1099番地 株

式会社日立製作所システム開発研究所内

(72) 発明者 山本 彰

神奈川県川崎市麻生区王禅寺1099番地 株

式会社日立製作所システム開発研究所内

(72) 発明者 佐藤 孝夫

神奈川県川崎市麻生区王禅寺1099番地 株

式会社日立製作所システム開発研究所内

(74) 代理人 弁理士 小川 勝男

最終頁に続く

(54) 【発明の名称】 ディスクアレイ制御方法

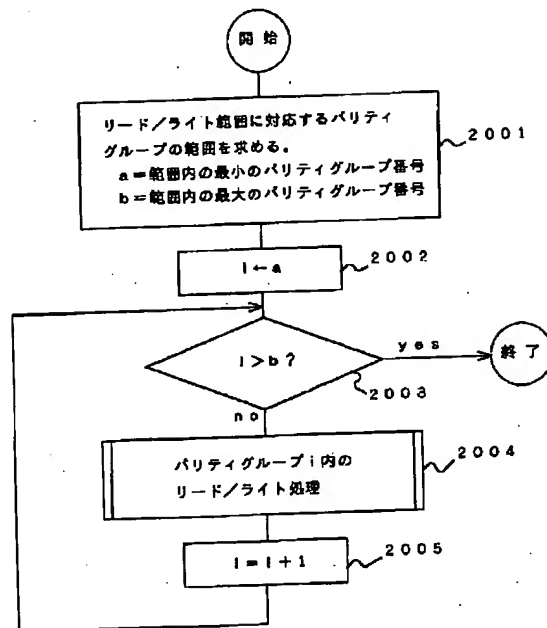
(57) 【要約】

【目的】 冗長構成を採用したディスクアレイにおけるディスク故障時の性能低下を防止する。

【構成】 ディスクアレイを構成するディスクに故障が生じた際、故障ディスクに格納されていたデータを上書きする上書き用の冗長データを選択し、故障したディスク内のデータまたは冗長データを上書き用の冗長データの上に上書きし、ディスク構成を管理する管理テーブルを更新する。

【効果】 ディスク故障時のディスクアクセス回数を低減することができ、ディスク故障時のディスクアレイの性能低下を防止することができる。

図 19



【特許請求の範囲】

【請求項1】データを格納するための少なくとも1のディスク装置と冗長データを格納するための少なくとも1のディスク装置とを備えたディスクアレイシステムにおけるディスクアレイ制御方法において、前記データを格納するディスク装置および冗長データを格納するディスク装置の故障を検知し、前記故障が検知された場合、前記冗長データを格納するディスク装置から任意のディスク装置を選択し、前記故障が検知された以外のディスク装置に格納されたデータおよび冗長データを用い前記故障が検知されたディスク装置のデータを復元し、該復元されたデータを前記選択により選択されたディスクに記録することを特徴とするディスクアレイ制御方法。

【請求項2】前記復元するステップと、前記記録するステップは、ホストコンピュータからのディスクアレイへのリードおよびライトアクセスとは非同期なバックグラウンド処理として行なわれることを特徴とする請求項1記載のディスクアレイ制御方法。

【請求項3】前記復元するステップと、前記記録するステップは、ホストコンピュータからの故障ディスクへのリードおよび/またはライトアクセス時に同期して行なわれることを特徴とする請求項2記載のディスクアレイ制御方法。

【請求項4】請求項1記載のディスクアレイ制御方法において、前記ディスクアレイシステムは、故障ディスク装置を代替する少なくとも1台のホットスタンバイ用の予備ディスクを備え、前記故障を検知するステップにおいて故障が検知された際、前記呼びディスクが全て代替用として用いられている場合に前記選択するステップ、復元するステップ、および記録するステップを実行することを特徴とする請求項1記載のディスクアレイ制御方法。

【請求項5】前記各ステップを繰り返し実施することを特徴とする請求項1記載のディスクアレイ制御方法。

【請求項6】多重度 n (≥ 1)のデータに対して m (≥ 1)個の冗長データを記録する、ディスクの冗長構成を採用したディスクアレイを制御するディスクアレイ制御方法であって、 k ($m \geq k \geq 1$)台のディスクが故障した時に、前記 m 個の冗長データの中から k 個の冗長データを選択するステップと、故障した k 台のディスク内に記録されていた k 個のデータを前記ステップで選択した k 個の冗長データの格納されていた領域に書き込むステップを有し、前記故障した k 個のディスク内に記録されていた k 個のデータを正常なディスク内の冗長データの格納されている k 個の領域に書き込むことによりディスクの冗長構成を再構成することを特徴とするディスクアレイ制御方法。

【請求項7】前記ディスクの冗長構成の再構成処理が、ホストコンピュータから前記ディスクアレイへのリード/ライトアクセスとは非同期なバックグラウンド処理と

して行なわれることを特徴とする請求項6記載のディスクアレイ制御方法。

【請求項8】前記ディスクの冗長構成の再構成処理が、ホストコンピュータから故障ディスクへのリード/ライトアクセス時に、該リード/ライトアクセスと同期して行なわれることを特徴とする請求項6記載のディスクアレイ制御方法。

【請求項9】前記ディスクアレイは、さらに s (≥ 1)台のホットスタンバイの予備ディスクを有し、ディスク故障時に、全ての予備ディスクを使いきって、予備ディスクがなくなった状態で、更にディスクが故障したときに、前記ディスクの冗長構成の再構成処理を行うことを特徴とする請求項6記載のディスクアレイ制御方法。

【請求項10】前記ディスクアレイは、さらに s (≥ 1)台のホットスタンバイの予備ディスクを有し、 k ($m + s \geq k > s$)台のディスクが故障した時に、故障ディスク内のデータおよび冗長データの予備ディスクへの回復と前記ディスクの冗長構成の再構成とを行うことを特徴とする請求項6記載のディスクアレイ制御方法。

【請求項11】前記冗長データが2以上であって、前記ディスクの冗長構成の再構成を複数回繰り返すことを特徴とする請求項6記載のディスクアレイ制御方法。

【発明の詳細な説明】

【0001】

【産業上の利用分野】本発明は、コンピュータシステムにおける外部記憶装置の制御方法にかかり、特に、複数のディスク装置により冗長構成をなすディスクアレイ装置の障害時の性能低下を防止を考慮したディスクアレイ制御方法に関する。

【0002】

【従来の技術】近年、コンピュータシステムにおけるI/Oネックを解決する為の、一つのアプローチとして、冗長構成のディスクアレイが注目され、技術開発が活発に行われている。冗長構成のディスクアレイはRAID (Redundant Array of Inexpensive Disks) と呼ばれ、その構成方式として、RAID3、RAID4、RAID5が知られている。RAIDに関しては、例えば、デー・バターソン、ジー・ギブソン、アール・カツ等による、「ア ケース フォー リダンダントアレイズ オブ インイクスペンシブ ディスクス (レイド)、エーシーエムシグモッド カンファレンス プロシーディングス (D.Patterson, G.Gibson, R.Katz, "A Case for Redundant Arrays of Inexpensive Disks (RAID)", ACM SIGMOD conference proceedings, 1988, pp. 109-116)」に述べられている。

【0003】RAID3とは、ディスクアレイを構成するディスクをデータを格納するディスク (以下、データディスクと呼ぶ) と冗長データを格納するディスク (以下、パリティディスクと呼ぶ) とに分け、ホストコンピュータから見える論理的なディスクのアドレスが、実際

3

にデータを格納する物理ディスク間を切り替わる大きさ(以下、データストライプサイズと呼ぶ)を、1ビットないし1バイト、あるいは1セクタ等と小さくとしたものである。これにより、典型的なディスクアクセスが、アレイを構成する全てのディスクにまたがるようにしてある。RAID3は、大量データの一括転送能力が優れている。

【0004】一方、RAID4とは、RAID3と同様に、ディスクアレイを構成するディスクをデータディスクとパリティディスクを分け、データストライプサイズを典型的なディスクアクセスのデータサイズ以上に10して、典型的なディスクアクセスが、1台の物理ディスクに収まるようにしたものである。RAID4では、複数個のライトアクセス間で、パリティディスクが競合するという問題がある。

【0005】また、RAID5は、RAID4におけるパリティディスク競合の問題を解決するために、冗長データをディスクアレイを構成する複数台のディスク間に分散して格納するようにしたものである。RAID5において、あるディスクが故障したときの、デグレード時の動作や性能、予備のホットスタンバイディスクへの故障ディスクのデータの回復の方法や性能に関しては、例えば、エム・ホランド、ジー・ギブソン、等による「パリティ デクラスタリング フォー コンティニアス オペレーション イン リダンダント ディスク アレイズ、エーエスピーロスーバイ プロシーディングズ」(M.Holland, G.Gibson, "Parity Declustering for Continuous Operation in Redundant Disk Arrays", ASPLOS-V Proceedings, October 12-15, 1992, pp.2-9.)に開示されている。

【0006】

【発明が解決しようとする課題】上記従来技術によれば、ディスクアレイ内のディスク故障に対する信頼性は向上されるが、ディスク故障が存在するときに、データのリード/ライトを行なう場合、ディスク故障のない時に比べ、ディスクアクセス回数が増大するという問題があった。

【0007】例えば、RAID4のディスクアレイでは、通常、少量のデータの読み出しであれば、ディスクのアクセス回数は1回で済む。しかし、ディスクアレイ内のあるデータディスクが故障したとき、故障したディスク内のデータを読み出そうとすると、正常なデータディスク、及び冗長ディスクから、リードしようとするデータを復元するのに必要なデータおよび冗長データを読み出し、排他的論理和を取るなどして、目的とするデータを復元する必要がある。従って、データの冗長構成の単位であるパリティグループが、 n 個のデータと m 個の冗長データにより構成されるものを考えた場合、少なくとも、残りの正常なデータディスク $n-1$ 台と、冗長ディスク1台から、データを読み出さなければならないこ

4

とになる。このため、最低でも計 n 台のディスクをアクセスする必要が生じる。つまり、故障したディスク内のデータを読み出そうとした場合には、ディスクアクセス回数が、最低でも通常のリード時の n 倍に増加することになる。

【0008】同様に、RAID5においても、通常、少量のデータの読み出し時には、ディスクのアクセス回数は1回で済む。これに対して、ディスクアレイ内のディスクが故障したとき、故障ディスク内のデータを読み出すためには、データの復元を行うために、残りの正常なディスクのうち、データを格納しているディスク $n-1$ 台と、冗長データを格納しているディスク少なくとも1台をアクセスして、データの復元に必要な、 $n-1$ 個のデータと、少なくとも1個の冗長データを、読み出す必要がある。従って、この場合にも、最低 n 台のディスクをアクセスする必要がある。

【0009】以上、データの読み出しの場合を例に説明したが、同様のことがデータの書き込み時にも問題となる。

【0010】従って、本発明の目的は、ディスクアレイに関する上述したような問題を解決し、ディスク故障時における故障ディスク内のデータに対するリード/ライト時の、ディスクアクセス回数の増加を防止したディスクアレイの制御方法を提供することにある。

【0011】

【課題を解決するための手段】上記目的を達成するために、本発明におけるディスクアレイの制御方法は、データを格納するための少なくとも1のディスク装置と冗長データを格納するための少なくとも1のディスク装置とを備えたディスクアレイシステムにおけるディスクアレイ制御方法において、前記データを格納するディスク装置および冗長データを格納するディスク装置の故障を検知し、前記故障が検知された場合、前記冗長データを格納するディスク装置から任意のディスク装置を選択し、前記故障が検知された以外のディスク装置に格納されたデータおよび冗長データを用い前記故障が検知されたディスク装置のデータを復元し、該復元されたデータを前記選択により選択されたディスクに記録することを特徴としている。

【0012】より詳しくは、多重度 n (≥ 1) のデータに対して m (≥ 1) 個の冗長データを記録する、ディスクの冗長構成を採用したディスクアレイにおいて、 k ($m \geq k \geq 1$) 台のディスクが故障した時に、まず、 m 個の冗長データのうち、 k 個の冗長データを選択し、故障した k 台のディスク内に記録されていた k 個のデータまたは冗長データのうち、データと、前記ステップ1で選択しなかった冗長データを復元し、正常なディスク内の、前記ステップ1で選択した k 個の冗長データの格納されていた領域に書き込む。

【0013】

5

【作用】本発明によれば、ディスクアレイを構成するディスクの一部に故障が生じた際、冗長データを記録したディスクから故障が生じたディスクに格納されたデータまたは冗長データを重ね書きするディスクを選択し、この選択されたディスクに故障により失われたデータを復元して記録する。このため、ホストコンピュータからのリード/ライト要求に対し、ディスク故障が存在する場合であっても復元されたデータを用いて処理を行なうことができ、アクセス回数の増大を抑えることが可能である。従って、ディスク故障時における性能低下を防止することができる。

【0014】

【実施例】以下、本発明につき、図面を用い詳細に説明する。

【0015】図1は、本発明が適用される計算機システムの構成図であり、計算機システム内におけるディスクアレイの位置付けを簡略に示したものである。図1において、100はディスクアレイであり、ディスクアレイ内においてデータのリード/ライト等の制御を行なうディスクアレイコントローラ200と、複数のディスクドライブ300から構成されている。400は、ディスクアレイ100が接続されるホストコンピュータである。このような計算機システムにおいて、ディスクアレイ100は、複数台のディスクドライブ300をホストコンピュータ400に論理的に一つのディスク装置に見せる機能を持つ。以下、本実施例では、RAID5構成のディスクアレイについて述べるが、本発明は、RAID4構成のディスクアレイにも適用できることはいうまでもない。

【0016】図2は、ホストコンピュータから見たディスクアレイのイメージである論理ディスクと、実際にデータを格納する物理ディスクとの間の対応を示したマッピング機能図である。

【0017】図2において、200はディスクアレイコントローラ、301はホストコンピュータから見えるディスクアレイのイメージである論理ディスク、300は実際にデータを格納するディスクドライブ、即ち物理ディスクである。501は論理ディスクにおけるデータの管理単位であるデータブロック、500、600は物理ディスクにおけるデータの管理単位であり、それぞれデータを格納するデータブロック、冗長データを格納するパリティブロックである。図2は、1台の論理ディスクが $n+m$ 台の物理ディスクによって実現される構成を示している。 $n+m$ 台の物理ディスクを、DISK#1, DISK#2, ..., DISK# n , DISK# $n+1$, ..., DISK# $n+m$ と表現する。ここで、論理ディスク301のデータブロック501、物理ディスクのデータブロック500、パリティブロック600は全て同一の大きさであるとする。この大きさは、論理ディスク301内のデータが、物理ディスク300間を切り替

6

わるデータの折返し単位（ストライプサイズと呼ぶ）であるとともに、冗長データが、物理ディスク間を切り替わる冗長データのストライプサイズでもある。

【0018】図2に示したアレイ構成は、RAID5のレフトシンメトリック（left symmetric）と呼ばれる構成である。冗長データは n 個のデータに対して、 m 個作成され、 n 多重のデータに対して、冗長度が m の冗長構成となる。700はデータの冗長構成の単位であるパリティグループであり、 n 個のデータブロックと、それに対応した m 個のパリティブロックからなる。論理ディスク301と物理ディスク300の対応は次のようになる。論理ディスク301内のデータを、先頭から n 個の連続したデータブロック毎に区切る。この n 個のデータブロック毎に、 m 個の冗長データを作成し、 n 個のデータと m 個の冗長データで一つのパリティグループ700を構成する。各パリティグループ700には、論理ディスク内のアドレスの順に、パリティグループ番号を付ける。また、 m 個の冗長データには、一連の記号P1, P2, ..., P m を付ける。

【0019】パリティグループ700内のデータと冗長データの、物理ディスクへの格納は次のように行なわれる。

【0020】論理ディスク301内のデータブロックに対して、その番号を n の剰余系で考えて、1～ n の整数値を与える。次に、対応する m 個の冗長データP1, P2, ..., P m に、それぞれ整数値 $n+1$, $n+2$, ..., $n+m$ を与え、パリティグループ内の、 n 個のデータブロック500と m 個のパリティブロック600に対して、1から $n+m$ までの整数値を与える。次に、パリティグループ番号1のデータブロック500およびパリティブロック600を、それらに与えられた整数値1と等しいディスク番号を有する物理ディスク300に格納する。つまり、パリティグループ番号1のパリティグループについては、データブロック1を、DISK1に格納し、データブロック2, 3, ..., n をそれぞれ、DISK2, DISK3, ..., DISK n に格納する。パリティブロックについては、P1からP m までを、それぞれDISK $n+1$ からDISK $n+m$ に格納することになる。パリティグループ番号 j （但し、 $j \geq 2$ ）のパリティグループについては、前述の方法によって整数値が与えられた各データブロック500及びパリティブロック600を、パリティグループ番号 $j-1$ 内の同じ整数値が与えられたデータブロック500およびパリティブロック600が格納される物理ディスクのディスク番号を、1だけ減じたディスク番号を持つ物理ディスク300に格納する。即ち、図2に示すようにパリティグループ番号が1増す毎に、各ブロックの格納位置が、左へ1つづつずれてゆくことになる。

【0021】図3には、ディスクアレイの構成図を示す。図3において、200はディスクアレイコントローラ

7
 ラ、300はディスクドライブ、400はホストコンピュータである。ディスクアレイコントローラ200は、CPU210、冗長データの生成、及び障害ディスク内のデータを復元するパリティ生成・データ復元回路220、ホストコンピュータ400との間でデータの受渡し等を行なうホストインタフェース230、それぞれ対応するディスクドライブを制御するドライブコントローラ240、メモリ260、及びメモリ260へのデータのリード/ライトを制御するメモリコントローラ250を含んで構成される。これらディスクアレイコントローラ200の構成要素は、バス270により互いに接続されている。

【0022】メモリ260には、各種の処理をCPU210が実施するためのプログラムと、各ディスクドライブ300の状態を管理するディスク状態管理テーブル261、冗長データを管理するためのパリティ管理テーブル262、ディスク故障からの回復処理の状態を示す回復処理管理テーブル263、各パリティグループ毎にデータ及び冗長データがどのディスクドライブに格納されているかを示すパリティ構成管理テーブル264、および回復処理に用いられ、データが上書きされている冗長データの種別を記録する上書き用パリティ管理テーブル265等の制御テーブルが格納されている。

【0023】以下、本実施例において中心的役割を果たすディスク故障時における冗長構成の再構成の方法について説明する。なお、ここでは、説明の便宜上、6台のディスクで構成され、各パリティグループが4個のデータブロックと2個のパリティブロックからなるディスクアレイを想定して説明する。

【0024】図4、図5に、ディスク状態管理テーブル261の論理的な構成を示す。ディスク状態管理テーブル261には、各ディスク毎に、それが正常であるか、故障中であるかが示されている。図4は、全てのディスクが正常状態にある場合のディスク状態管理テーブルを示しており、図5はDISK#2のディスクが故障している場合のディスク状態管理テーブルを示している。

【0025】図6、図7にはパリティ管理テーブル262の論理的な構成を示す。パリティ管理テーブルは、個々の冗長データがそれぞれ冗長データとして使用中であるか、ディスク故障により故障ディスクのデータを上書きされたものであるか（使用中）を管理する。図6に示すパリティ管理テーブルは、2つの冗長データがともに冗長データとして使用されていることを示しており、図7に示すパリティ管理テーブルは、冗長データP2がデータの上書き用に用いられていることを示している。

【0026】図8に、回復処理管理テーブル263の論理的構成を示す。回復処理管理テーブル263は、パリティグループ毎に、故障したディスク内のデータを冗長データが記録されているディスクに上書きすることにより回復済みか、まだ未回復であるかを管理する。

【0027】図9乃至図14には、パリティ構成管理テーブル264の論理的構成を示す。パリティ構成管理テーブル264は、各パリティグループ毎に、各物理ディスク内に格納されているデータ、冗長データを識別する情報が記録される物理論理変換部と、データおよび冗長データが、どの物理ディスクに格納されているかを記録する論理物理変換部とを有している。本実施例ではデータD1、D2、D3、D4をそれぞれ識別する情報として整数値1、2、3、4を、また、冗長データP1、P2を識別する情報として、整数値5、6を用いている。なお、図9から図14に示される各パリティ管理テーブルの意味するところは、以下に説明する本実施例の動作とともに説明する。

【0028】図16は、本実施例におけるディスク故障モニタ処理の流れを示すフローチャートである。ディスク故障モニタ処理では、ディスク故障の検知、回復処理の起動などを行う。

【0029】ディスク故障モニタ処理では、各ディスクの状態を調べ、ディスク故障の検知を行い（ステップ1701）、ディスクが故障していなければ、そのまま一定時間待ち、再びステップ1701へ戻ってディスク故障の検知を行なう（ステップ1708）。ステップ1701で、ディスクの故障が検知されると、故障ディスクの特定を行ない（ステップ1702）、以降、故障ディスクを制御から切り離して動作するために、制御テーブルの設定を行う（ステップ1703）。ここでは制御テーブルの設定として、ディスク状態管理テーブル261の設定とパリティ構成管理テーブル264の設定を行う。

【0030】ディスクアレイを構成する6台のディスクが全て正常である場合、ディスク状態管理テーブル261、および、パリティ構成管理テーブル264は、それぞれ図4、図9に示す状態にある。ここで、ステップ1701においてディスク故障が検知され、ステップ1702において故障を生じたディスクがDISK#2であることが特定されたとする。ステップ1703では、まず、図5に示すようにディスク管理テーブル261のDISK#2の状態を、正常を示す状態から故障を示す状態に変更する。次に、図9に示す正常時のパリティ構成管理テーブル264から、ディスク故障時のリード/ライト制御に使用する、故障時のパリティ構成管理テーブル261を作成する。具体的には、正常時のパリティ構成管理テーブルの物理論理変換部の物理ディスク番号2の列内のデータを無効とする。更に、論理物理変換部の、物理ディスク番号2のデータを格納していたエンタリ内のデータを無効にする。このようにして作成したパリティ構成管理テーブルを図10に示す。

【0031】制御テーブルの設定が終了すると、ディスクが故障する前に使用していた冗長データの個数、即ち、ディスク故障検知前の冗長度が1以上であったかど

うかチェックする(ステップ1704)。このチェックは、パリティ管理テーブル262を調べることににより行なうことができる。パリティ管理テーブル262内に、使用中の冗長データが u 個あれば、冗長度 u で動作していたことになる。上長度検査の結果、ディスク故障検知前に冗長度0で動作していた場合は、ディスク故障により、以降の動作が不可能となるので、ホストコンピュータに障害を報告し、停止する(ステップ1709)。ディスク故障検知前に冗長度1以上で動作していた場合は、ステップ1705へ進み、故障ディスクのデータを上書きする冗長データを、パリティ管理テーブル262から選択し、パリティ管理テーブル262にを更新する。ここでは、2つの冗長データのうち、2番目の冗長データであるP2を故障ディスクのデータ上書き用に選択し、パリティ管理テーブルのP2の状態情報を使用中から上書き用に変更する。変更後のパリティ管理テーブル262を図7に示す。また、選択した冗長データの種別、P2を上書き用パリティ管理テーブル265に記録する。

【0032】次に、回復処理管理テーブル263の、全てのパリティグループ番号の回復状態管理情報を未回復に設定する。また、回復済みのパリティグループに対するリード/ライト処理を制御するために、図11に示すような、データ上書き後のパリティ構成管理テーブルを作成する(ステップ1706)。それから、非同期の回復処理を起動し(ステップ1707)、回復処理終了後、ステップ1701へ戻り、ディスク故障のモニタリングを続ける。

【0033】次に、ステップ1707で実施される回復処理を説明する。

【0034】図17は回復処理の流れを示すフローチャートである。

【0035】図17において、 k は処理中のパリティグループの番号を保持する内部変数である。回復処理は、ディスクアレイに対するホストコンピュータからのリード/ライトとは非同期に行なわれる。また、回復処理は、パリティグループ番号1からパリティグループ番号順に行なわれる。

【0036】回復処理では、まず、内部変数 k の値を1にする(ステップ1801)。次に、内部変数 k が、ディスクアレイ内の最大のパリティグループ番号の値、 P_{GNmax} より大きいのかを判定する(ステップ1802)。内部変数 k が P_{GNmax} より大であれば、全てのパリティグループについて回復処理を完了したことになるので、処理を終了する。内部変数 k が P_{GNmax} 以下であれば、パリティグループ番号 k のパリティグループについて回復処理を行う(ステップ1803)。パリティグループ番号 k のパリティグループについて回復処理終了後、内部変数 k に1を加え、次のパリティグループの処理へ進む(ステップ1804)。

【0037】図18は、図17に示す回復処理1803の詳細な流れを示すフローチャートである。

【0038】まず、処理対象のパリティグループ(パリティグループ番号 k)が、回復済みか否かを、回復処理管理テーブル263により調べ、処理対象のパリティグループ回復済みであれば、回復処理を終了する(ステップ1901)。

【0039】一方、処理対象のパリティグループが未回復であれば、そのパリティグループ内で、故障ディスクに格納されていたデータ/冗長データの種別 c を求める(ステップ1902)。これは次の方法によって可能である。まず、処理対象のパリティグループ番号の値を、データブロック及びパリティブロックの総数 $n+m$ (本実施例では6)で除してやり、余りが0であれば $n+m$ を、余りが0以外であればその余りを、そのパリティグループのパリティレイアウト番号として求める。ここで、パリティレイアウトとは、パリティグループにおいて、各物理ディスクに、どのようにデータおよび冗長データを記録させるかという、データ及び冗長データの物理ディスクへの配置パターンをいう。多重度 n のデータに、 m 個の冗長データを付加する冗長構成のRAID5では、パリティレイアウトは、計 $(n+m)$ 個存在する。次に、図9に示す正常時のパリティ構成管理テーブル264を用い、該当するパリティレイアウト番号の列と、物理論理変換部の故障ディスクの物理ディスク番号の行で特定されるエントリ内の値を求める。これが、データ/冗長データの種別 c である。

【0040】次にステップ1902で求めた種別 c が、上書きに用いられる冗長データか否かを上書き用パリティ管理テーブル265を参照して判定する(ステップ1903)。種別 c が、上書きに用いられる冗長データであれば、上書き処理は不要であり、ステップ1908へジャンプする。種別 c が上書きに用いられる冗長データでなければ、ステップ1904へ進み、種別 c がデータか冗長データかを判定する。これは、ステップ1903同様、正常時のパリティ構成管理テーブル264を参照することにより可能である。種別 c がデータであれば、そのデータを回復するために、同一のパリティグループに属するデータと冗長データを格納しているディスクの番号を、パリティ構成管理テーブルから求め、それらを該当するディスクから読み出し(ステップ1905)、故障したディスク内のデータを回復する(ステップ1906)。種別 c が冗長データであれば、その冗長データを生成するため、パリティグループ内のデータを、ディスクから読み出し(ステップ1909)、それらを基に冗長データを生成する(1910)。

【0041】最後に、回復したデータまたは生成した冗長データを、上書き用に選択した冗長データの上に上書きし(ステップ1907)、回復状態管理テーブル263の処理対象のパリティグループのエントリに対応する

回復状態管理情報を未回復から回復済みへ変更して処理を終了する(ステップ1908)。

【0042】図19は、回復処理中の、ホストコンピュータから要求に対するリード/ライト処理の流れを示すフローチャートである。

【0043】ディスクアレイコントローラ200は、ホストコンピュータからのリード/ライト要求を受けると、まず、リード/ライト要求のあったデータの範囲に対応するパリティグループの範囲(a-b)を求める(ステップ2001)。ここで、aをリード/ライト範囲に対応する最小のパリティグループ番号とし、bを同様に最大のパリティグループ番号とする。パリティグループ番号a、bは、リード/ライト範囲の論理アドレスの最小値と最大値をそれぞれ、パリティグループ内のデータの多重度nとストライプサイズsとの積snで割った商として求められる。

【0044】次に、リード/ライト処理中のパリティグループの番号を保持する内部変数1に値aを代入する(ステップ2002)。内部変数1の値がbより大か否かを判定し(ステップ2003)。内部変数1がbより大であれば、リード/ライト範囲に対応するパリティグループの処理が完了しているので、処理を終了する。内部変数1がb以下であれば、パリティグループ1内のリード/ライト処理を行う(ステップ2004)。次に、内部変数1に1加算し、ステップ2003へジャンプして処理を続行する(ステップ2005)。

【0045】図20、および21に、ステップ2004の各パリティグループ毎のリード/ライト処理の流れを示す。

【0046】パリティグループ毎のリード/ライト処理では、まず、処理対象のパリティグループが、ディスク故障から回復済みか否かを、回復処理管理テーブル263を用いて判定する(ステップ2101)。回復済みであれば、回復済みのリード/ライト処理を行って処理を終了する(ステップ2107)。なお、回復済みのリード/ライト処理は、図11に示すようなデータ上書き後のパリティ構成管理テーブルを参照して、リード/ライトの対象となるデータの格納されたディスクを調べ、そのディスクに対するアクセスを行なうことにより処理を行うことができる。この時のリード/ライト処理は、データ多重度nに対し、冗長データの数m-1個の、故障したディスクのない、ディスクアレイに対するリード/ライト処理に帰着するので、詳しい説明は省略する。

【0047】処理対象のパリティグループが、未回復である時は、ホストコンピュータから要求された処理がリード処理であるかライト処理であるかにより分岐して処理を行う(ステップ2102)。

【0048】ホストコンピュータからの処理の要求がリード処理であるときには、まず処理対象のパリティグループ内のリード領域が、どの物理ディスクに対応するか

を求める(ステップ2103)。次に、リード領域に故障ディスクを含むか否かチェックする(ステップ2104)。リード領域に故障ディスクを含まない場合、ステップ2103で求めた、物理ディスクから、データを読み出し(ステップ2105)、データをホストコンピュータに転送して処理を終了する(ステップ2106)。

【0049】ステップ2104で、リード領域に故障ディスクを含むことがわかったときは、パリティグループ内の、故障ディスクに格納されるデータ以外のデータ全てと、故障ディスク内のデータを回復するのに必要な冗長データの一つ、正常なディスクから読み出す(ステップ2108)。そして、故障ディスクのデータを回復してステップ2106へ進み(ステップ2109)、先と同様に、リード対象のデータをホストコンピュータへ転送して処理を終了する。ここで、ステップ2107において読み出すデータおよび冗長データが、どの物理ディスクに格納されているかは、正常時のパリティ構成管理テーブル264から求めることができる。

【0050】ステップ2102で、ホストコンピュータからの処理の要求がライト処理と判定された場合には、図21のラベルWへジャンプし、処理を続ける。以下、図21を参照してライト処理につき説明する。

【0051】まず、正常時のパリティ構成管理テーブル264を用いてパリティグループ内のライト領域が、どの物理ディスクに対応するかを求める(ステップ2201)。次に、ライト領域に故障ディスクを含むか否かを判定する(ステップ2202)。ライト領域に故障ディスクを含まない場合、ステップ2201で求めた物理ディスクから、ライト領域に対応する更新前のデータ(旧データ)、および更新前の冗長データを読み出す(ステップ2203)。次に、ホストコンピュータから受け取ったライトデータと、ディスクからリードした、旧データ、および更新前の冗長データから、新しい冗長データを生成する(ステップ2204)。次に、新しい冗長データと、ライトデータを、それぞれ、該当する物理ディスクに書き込み処理を終了する(ステップ2205)。

【0052】以上により、処理対象のパリティグループ内のライトデータのディスクへの書き込みと、冗長データの更新が完了する。

【0053】ステップ2202で、ライト領域に故障ディスクを含んでいた場合、処理対象のパリティグループ内の物理ディスクでデータを格納しているもののうち、ライト領域に対応する物理ディスクを除いた、残りの物理ディスクからデータをリードする(2206)。次に、ライトデータと、ステップ2206でリードしたデータから、新しい冗長データを生成する(2207)。それから、故障ディスクを除いて、ライトデータと、ステップ2207で生成した新しい冗長データを物理ディスクにライトして処理を終了する(ステップ2208)。

【0054】以上説明したように、本実施例によれば、回復処理中にも、ホストコンピュータからのリード/ライトアクセスに対する処理を、回復処理と並行して行うことができる。なお、以上の説明では、6台のディスクによりデータの多重度4、冗長度2の冗長構成を採っているときに、2番のディスク(DISK#2)が故障したとき例に説明したが、その後さらに、他のディスクが故障した場合にも同様の処理手順により回復処理を行なうことができる。上述したDISK#2の故障に対する回復処理完了後に、更にDISK#4が故障した場合を例に、このとき用いられるパリティ構成管理テーブルの例を図12、図13、図14に示す。DISK#2の故障に対する回復処理完了後に、更にDISK#4が故障したときは、残りの冗長データP1にDISK#4に格納されていたデータを上書きすることにより、これを回復することが可能である。

【0055】以上述べたように、本実施例によれば、ディスク故障時に、ディスクアレイの冗長構成を、冗長度を1だけ落とした、故障ディスクの無い冗長構成に再構成するため、故障ディスク内のデータに対するリード/ライト処理におけるディスクアクセス回数を従来に比べ、低減することができる。また、冗長度が1だけ減るので、ライト時に更新すべき冗長データの数1だけ少なくなり、ライト時のライトペナルティを減らすこともできる。

【0056】また、本実施例におけるディスクアレイの制御方法は、RAID3構成のディスクアレイに適用することも可能である。例えば、多重度nのデータに1個の冗長データを付加する冗長構成のRAID3に本実施例における制御方法を適用した場合を考えてみる。データを格納するn台のディスクのうち、1台が故障したとすると、このとき、データのリード処理は、データをn-1台のデータディスクから、また、冗長データをそれを格納する1台のパリティディスクからそれぞれ読み出し、故障したディスクのデータを復元してそのデータをホストコンピュータに送ることにより行なわれる。また、データのライト処理は、ホストコンピュータから送られてきたデータから、冗長データを作成し、故障していないn-1台のデータディスクと、1台のパリティディスクに、それぞれデータと冗長データを書き込むことにより行なわれる。この場合、リード処理時には故障したディスク内のデータの復元が不要となり、ライト時には冗長データの生成が不要になるという利点がある。

【0057】以上説明した実施例では、故障したディスク内のデータが通常データ、および冗長データの何れの場合にもデータ回復を行なっているが、以下、第2の実施例として、故障したディスク内のデータが冗長データである場合にはその回復処理を行なわないものについて説明する。

【0058】本実施例においては、システム構成等は第

1の実施例と同じであり、その説明は省略し、ディスク故障検出時の処理のうち、第1の実施例とは異なる部分について説明する。

【0059】本実施例では、図17に示す回復処理のうち、ステップ1803におけるパリティグループ毎の回復処理が第1の実施例とは多少異なる。図22に本実施例によるパリティグループ毎の回復処理のフローチャートを示す。なお、図22では、図18と同様の処理が行なわれる部分については、図18と同一の符号を用いている。

【0060】第1の実施例においても説明したように、まず、処理対象となっているパリティグループが既に回復処理が行なわれたものかどうか判定し(ステップ1901)、回復済みであれば処理を終了し、未回復の場合、次のステップ1902に進み、処理対象のパリティグループのデータの種別を求める。処理対象のパリティグループ内のデータの種別が求まると、それに続いて、求めた種別が冗長データを示している場合には、そのままステップ1908に進んで回復状態管理情報の設定を行ない本処理を終了する。また、ステップ1902で求めた種別が冗長データでない場合には、ステップステップ1905へ進み、ステップ1906、およびステップ1907、を実行してデータを回復し、ステップ1908で回復状態管理情報の設定を行なって処理を終了する。

【0061】本実施例における回復処理完了後の、パリティ構成管理テーブルを図15に示す。ライト時の冗長データ更新、また、更に他のディスク故障が生じた場合のデータ回復には、図15に示す、パリティ構成管理テーブルを使用する。図16に示すように、本実施例では冗長データに対する回復処理を行なわないため、故障ディスク内に冗長データを記録していたパリティグループに対応するパリティレイアウトに対しては、上書きすべく選択した冗長データが残っている。即ち、パリティグループによって、保持される冗長データの種別が異なる場合が生じる。このため、ライト時の冗長データ更新や、更にディスクが故障した時のデータ復元処理が若干複雑になってしまう。しかし、故障ディスク内の冗長データについて回復処理を行なわないことにより、ディスクアレイの冗長構成の再構成に要する回復処理の時間、処理データ量を、削減できるというメリットを有している。第1及び第2の実施例では、故障ディスク内のデータ(または冗長データ)の回復処理をホストコンピュータからのディスクアクセスとは非同期にバックグラウンド処理として行なうものについて説明した。次に、第3の実施例として、ホストコンピュータから、故障したディスク内のデータに、リードまたはライトアクセスがあった時、要求データを対応する上書き用に選択された冗長データ上に復元する、オンデマンドでのデータ回復の方法について説明する。

15

【0062】図23、図24を用いて、ホストコンピュータからのリード/ライト要求の処理に同期して回復処理を行う場合の、パリティグループ毎のリード/ライト処理の方法を示す。なお、本実施例においても他の部分は第1の実施例と同様であるので、その説明は省略する。

【0063】ホストコンピュータからのリード/ライト要求があったときには、まず、処理対象のデータを含むパリティグループが回復済みか否かをチェックする(ステップ2301)。このパリティグループが回復済みである場合には、回復済みのリード/ライト処理を行う(ステップ2307)。ここでは、図20に示したリード/ライト処理におけるステップ2107と同様であり、説明は省略する。

【0064】処理対象のパリティグループが未回復であった場合、ホストコンピュータからの処理要求がリードかライトかによって異なる分岐先に進み処理を行う(ステップ2302)。ホストコンピュータからの処理要求がリードのときには、まず、処理対象のパリティグループ内のリード領域が格納される物理ディスクを求める(ステップ2303)。次に、求められた物理ディスクに故障ディスクを含むか否かをチェックする(ステップ2304)。故障ディスクを含まない場合には、ステップ2303で求めた物理ディスクからデータを読み出し(ステップ2305)、そのデータをホストコンピュータに転送して処理を終える。(ステップ2306)

ステップ2304で、リード領域に故障ディスクを含むことがわかったときは、そのパリティグループに属し、故障ディスク以外のディスクに格納されているデータ全てと、故障ディスク内のデータを回復するのに必要な冗長データを読み出し(ステップ2308)、それら読み出したデータを用いて故障ディスクに格納されていたデータを回復する(ステップ2309)。その後、ステップ2309で回復した故障ディスクのデータをディスク故障モニタで選択して、上書き用パリティ管理テーブルに記録されている冗長データの上に上書きする(ステップ2310)。さらに、回復処理管理テーブル内の処理対象のパリティグループに対応する回復状態管理情報に回復済みを設定し(ステップ2311)、ステップ2306へ進んでリード対象のデータをホストコンピュータへ転送する。

【0065】ステップ2302で、ホストコンピュータからの処理要求がライトと判定された場合には、図24のラベルWへジャンプして処理を続ける。

【0066】まず、図9に示す正常時のパリティ構成管理テーブルを用いて、処理対象のデータが属するパリティグループ内のライト領域が、どの物理ディスクに対応しているかを求める(ステップ2401)。次にライト領域に故障ディスクを含むか否かを判定する(ステップ2402)。ライト領域に故障ディスクを含まない場合に

16

は、ステップ2401で求めた物理ディスクから、ライト領域に対応する更新前のデータ(旧データ)、および更新前の冗長データを読み出す(ステップ2403)。次に、ホストコンピュータから受け取ったライトデータと、ディスクから読み出したライト処理前の旧データと、更新前の冗長データから新しい冗長データを生成する(ステップ2404)。次に、新しい冗長データとライトデータをそれぞれ物理ディスクに書き込んで処理を終了する(ステップ2405)。以上により、処理対象のパリティグループ内のライトデータのディスクへの書き込みと、パリティの更新が完了する。

【0067】一方、ステップ2402において、ライト領域に故障ディスクを含むことが分かった場合には、処理対象のデータが属するパリティグループ内の物理ディスクでデータを格納しているもののうち、ライト領域に対応する物理ディスクを除いた、残りの物理ディスクからデータを読み出す(ステップ2406)。次に、ライトデータと、ステップ2406で読み出したデータから新しい冗長データを生成する(ステップ2407)。それから、故障ディスクを除いた健全なディスクのライトデータをそれぞれのディスクに書き込み、ステップ2407で生成した新しい冗長データのうち、故障ディスクのデータを上書きする冗長データを除く冗長データをそれぞれ対応するディスクにライトし、故障ディスクへのライトデータをディスク故障モニタにより選択された上書き用の冗長データの上に書き込み(ステップ2408)、回復状態管理情報をセットして処理を終了する(ステップ2409)。以上、述べた処理方法により、ホストコンピュータからのリード/ライト要求の処理に同期した、故障ディスク内のデータの正常なディスク内の冗長データ上への書き込み処理(回復処理)を行うことができる。なお、本実施例によるオンデマンドでのデータ回復の方法では、故障ディスク内のデータのみを復元し、故障ディスク内の冗長データは復元していない。

【0068】以上説明した第3の実施例では、故障ディスク内のデータの回復処理をホストコンピュータからリード/ライト要求に同期して行なうため、ホストコンピュータからリード/ライト要求の処理に際して復元されたデータを用いることができ、バックグラウンドで回復処理を行なう場合に比べ、回復処理自体でのディスクアクセス回数を低減させることができる。また、本実施例による回復処理を、例えば、第1の実施例で説明したようなバックグラウンドの回復処理と並行して行なうこともできる。この場合には、バックグラウンドの回復処理時に、回復状態管理テーブルを、パリティグループ番号が小さいものから検索して、最もパリティグループ番号が小さい未回復のパリティグループを見つけて回復処理を行えばよい。

【0069】第3の実施例では、ホストコンピュータからのリード/ライト要求に同期して、故障ディスク内の

データにリードまたはライトが直接行われた場合にのみ故障ディスク内のアクセス対象となったデータを回復していたが、リード/ライトの対象となるデータが属するパリティグループ内の故障ディスクに格納されたデータ及び冗長データを上書きすべく選択された冗長データ上に回復してもよい。以下、第4の実施例として、このようなデータの回復方法について、図25及び図26を用い説明する。

【0070】第3の実施例と同様に、まず、ステップ2301で、処理対象のパリティグループが回復済みか否かをチェックし、回復済みである場合には、回復済みのリード/ライト処理を行う(ステップ2307)。処理対象のパリティグループが未回復であった場合には、ホストコンピュータからの処理要求がリードかライトかによって分岐して処理を行う(ステップ2302)。ホストコンピュータからの処理要求がリードの時には、まず処理対象のパリティグループ内のリード領域が、どの物理ディスクに対応するか(ここでは、物理ディスクd~eが該当するものとする)を求める(ステップ2303)。次に、当パリティグループの故障ディスク内のデータ及び冗長データの種別fを判定する(ステップ2601)。種別fが上書き用に選択された冗長データである場合は、物理ディスクd~eからデータをリードし(ステップ2602)、回復処理管理テーブルの回復状態管理情報を回復済みにセットする(2603)。さらに、データをホストコンピュータに転送して処理を終える(ステップ2306)。

【0071】ステップ2601で種別fが上書き用に選択されていない冗長データである場合には、当パリティグループ内の全てのデータをディスクから読み出し(ステップ2604)、故障ディスク内の冗長データを生成し(ステップ2605)、ステップ2605で生成した故障ディスク内の冗長データを上書き用に選択されている冗長データに上書きしてステップ2603へ進む(ステップ2606)。

【0072】また、ステップ2601で、種別fがデータであった場合は、当パリティグループ内の、故障していないディスクに含まれる全てのデータと、冗長データを一つディスクから読み出し(ステップ2607)、故障ディスク内のデータを回復した後(ステップ2608)、故障ディスク内のデータを上書き用として選択されている冗長データに上書きしてステップ2603へ進む(ステップ2609)。

【0073】ステップ2302で、ホストコンピュータからの処理要求がライトであるときには、図26のラベルWへジャンプする。

【0074】ステップ2401でパリティグループ内のライト領域を物理ディスク番号(先程と同様に、物理ディスクd~eが該当するものとする)に変換し、ステップ2701へ進み、当パリティグループ内の故障ディ

ク内のデータあるいは冗長データの種別(fとする)を判定する。種別fが上書き用に選択されている冗長データである場合、物理ディスクd~eから旧データを読み出し、更に故障していないディスクから当パリティグループの冗長データを読み出す(ステップ2702)。次に、ステップ2702で読み出した旧データ、旧冗長データと、ホストコンピュータからのライトデータとから新しい冗長データを生成する(ステップ2703)。そして、新しい冗長データとライトデータをディスクに書き込み(ステップ2704)、回復処理管理テーブルの回復状態管理情報に回復済みをセットして処理を終了する(ステップ2705)。

【0075】ステップ2701で種別fが上書き用に選択されていない冗長データである場合は、物理ディスクd~e以外のディスクから当パリティグループの残りのデータを読み出し(ステップ2706)、この読み出したデータとライトデータとから、新しい冗長データを生成し(ステップ2707)、ライトデータおよび健全なディスクに書き込まれる冗長データをそれぞれ該当するディスクへ書き込み、ステップ2705へ進む(ステップ2708)。

【0076】また、ステップ2701で種別fがデータである場合は、当パリティグループ内のデータを格納しているディスクのうちディスクd~eを除くディスクからデータをリードし、更に必要なら冗長データをも読み出し(ステップ2709)、ライトデータとステップ2709で読み出したデータ、冗長データから新しい冗長データを生成し、更に必要であれば、故障ディスク内のデータを復元する(ステップ2710)。そして、故障していない健全なディスクのライトデータをディスクライトし、上書きしない冗長データをディスクライトし、ライトデータまたはステップ2710で復元したデータとして得られる、故障ディスク内のデータを、上書き用として選択されている冗長データ上に上書きし(ステップ2711)、ステップ2705へ進む。

【0077】以上説明した処理により、ホストコンピュータからのリード/ライトに同期して、リード/ライト対象のデータが属するパリティグループ内のデータまたは冗長データであって、故障ディスクに格納されたデータまたは冗長データを上書き用として選択した冗長データ上に回復することができる。

【0078】以上、本発明について、4つの実施例を基に説明してきたが、以上説明した実施例によれば、RAID4、あるいはRAID5によるディスクアレイにおいて、一部のディスクに故障が生じた際、ディスクアレイ内の故障ディスクを排除して、健全なディスクのみでデータ構成を再編成するので、ディスク故障時におけるデータアクセス処理の性能低下を防止することができる。図27、28にそれぞれRAID4、RAID5における本発明と従来方式のアクセス回数につき、表にし

て示す。また、以上説明してきた実施例では、データ格納用のディスクと冗長データ格納用のディスクからなるディスクアレイについて説明したが、本発明は、さらに、ディスク故障時に故障ディスクの代替用として用いられるホットスタンバイディスクを有するディスクアレイに適用することもできる。ホットスタンバイディスクを有するディスクアレイでは、ディスク故障時に、故障ディスク上のデータや冗長データをホットスタンバイディスク上に復元する。この場合でも、ディスクの故障が重なり、全てのホットスタンバイディスクを代替として使いきっているような状態で、更にディスクの故障が発生したようなときに本発明が有効となる。

【0079】さらに、本発明をRAID3のディスクアレイに適用した場合には、ホストコンピュータからのデータリード要求に対する処理における故障ディスクのデータを復元する復元処理、あるいはホストコンピュータからのデータライト要求に対する処理における冗長データを生成する処理を省けるといった利点がある。

【0080】

【発明の効果】本発明によれば、ディスクアレイを構成する一部のディスクに故障が生じたときであっても、故障ディスクに対するリード/ライト時のディスクアクセス回数増加を防ぐことができる。また、このような場合における正常なディスクに対するライト時のディスクアクセス回数についても、低減することが可能となり、ディスク故障時の性能低下を抑えることができる。

【図面の簡単な説明】

【図1】本発明が適用される計算機システムの構成図である。

【図2】論理ディスクと物理ディスクの対応関係を示すマッピング機能図である。

【図3】ディスクアレイの構成図である。

【図4】正常時のディスク状態管理テーブルの論理的構成図である。

【図5】ディスク故障時のディスク状態管理テーブルの論理的構成図である。

【図6】正常時のパリティ管理テーブルの論理的構成図である。

【図7】ディスク故障時のパリティ管理テーブルの論理的構成図である。

【図8】回復処理管理テーブルの論理的構成図である。

【図9】正常時のパリティ構成管理テーブルの論理的構成図である。

【図10】ディスク故障時のパリティ構成管理テーブルの論理的構成図である。

【図11】再構成後のパリティ構成管理テーブルの論理的構成図である。

【図12】ディスク2故障時の回復処理完了後のパリティ構成管理テーブルの論理的構成図である。

【図13】ディスク2の回復処理完了後、更にディスク4が故障したときのパリティ構成管理テーブルの論理的構成図である。

【図14】ディスク2、4の回復処理完了後のパリティ構成管理テーブルの論理的構成図である。

【図15】本発明の第2の実施例における再構成後のパリティ構成管理テーブルの論理的構成図である。

【図16】ディスク故障モニタの処理フローチャートである。

【図17】非同期の回復処理のフローチャートである。

【図18】パリティグループ毎の回復処理のフローチャートである。

【図19】リード/ライト処理のフローチャートである。

【図20】パリティグループ毎のリード/ライト処理のフローチャートである。

【図21】パリティグループ毎のリード/ライト処理のフローチャートである。

【図22】本発明の第2の実施例によるパリティグループ毎の回復処理のフローチャートである。

【図23】リード/ライトに同期した回復処理を行う場合のパリティグループ毎のリード/ライト処理のフローチャートである。

【図24】リード/ライトに同期した回復処理を行う場合のパリティグループ毎のリード/ライト処理のフローチャートである。

【図25】パリティグループ毎のリード/ライト処理のフローチャートである。

【図26】パリティグループ毎のリード/ライト処理の続きである。

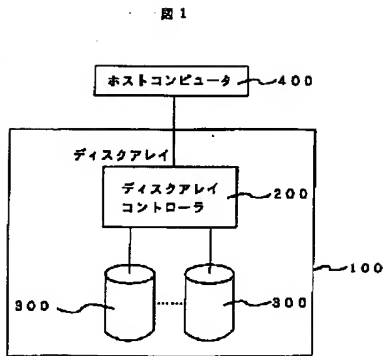
【図27】RAID4における本発明と従来技術との効果の差を示す比較図である。

【図28】RAID5における本発明の従来技術との効果の差を示す比較図である。

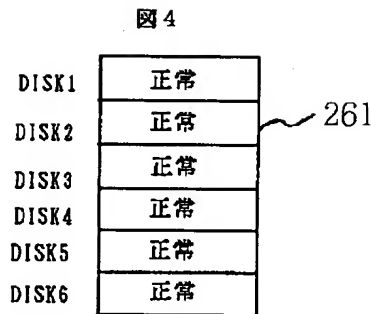
【符号の説明】

100…ディスクアレイ、200…ディスクアレイコントローラ、210…CPU、220…パリティ生成・データ復元回路、230…ホストインタフェース、240…ドライブコントローラ、250…メモリコントローラ、260…メモリ、270…バス、300…ディスクドライブ、301…論理ディスク、400…ホストコンピュータ、500、501…データブロック、600…パリティブロック、700…パリティグループ、701…パリティグループ番号である。

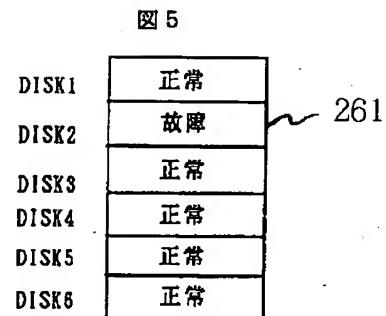
【図1】



【図4】



【図5】



【図2】

図2

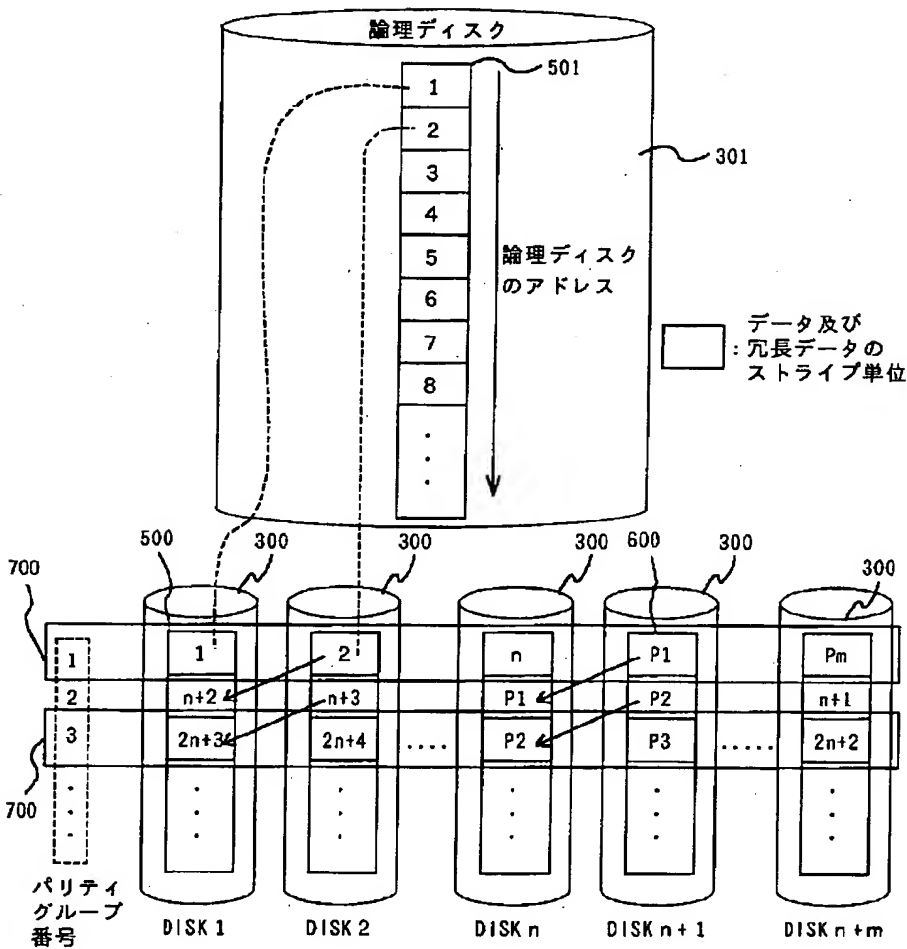
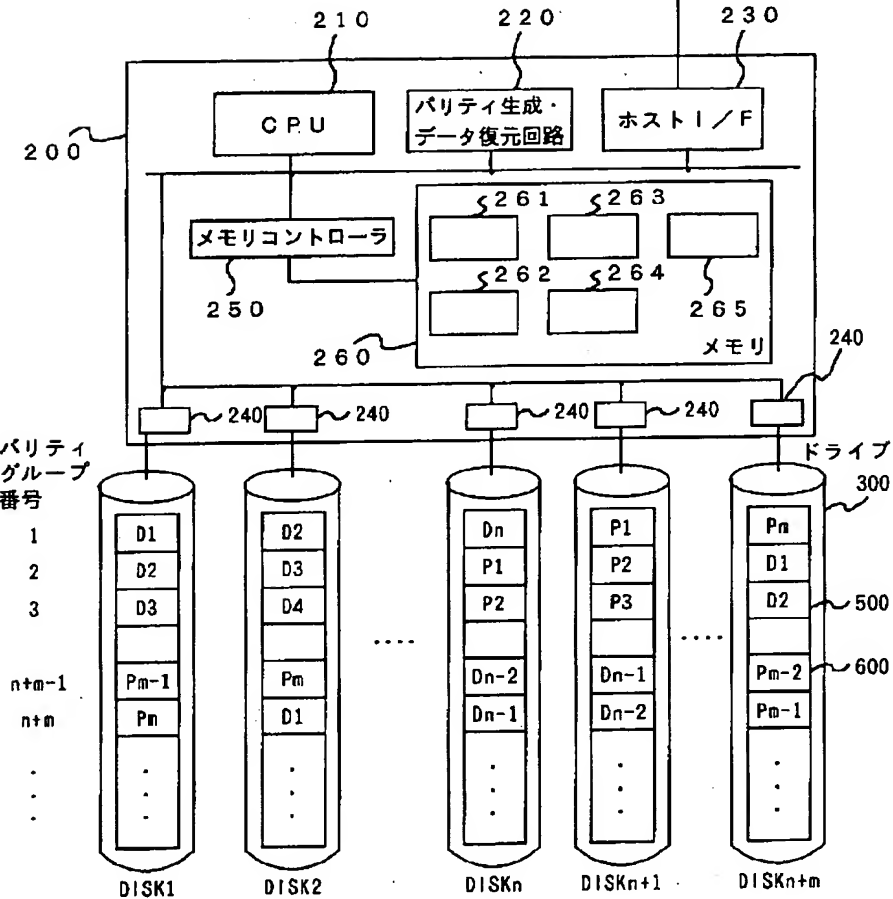


图 3



【图 7】

图 7

パリティ種別	状態情報
P1	使用中
P2	使用中

パーティ種別	状態情報
P1	使用中
P2	上書き用

【図8】

図 8

ハザード番号	回復状態管理情報
1	回復済み
2	回復済み
3	回復済み
4	回復済み
k	回復済み
k+1	未回復
PGN _{max}	未回復

263

【図9】

図 9

ハザード外番号	物理論理変換部	論理物理変換部
	物理ディスク番号	テープ/ハザード番号
	1 2 3 4 5 6	1 2 3 4 5 6
1	1 2 3 4 5 6	1 2 3 4 5 6
2	2 3 4 5 6 1	6 1 2 3 4 5
3	3 4 5 6 1 2	5 6 1 2 3 4
4	4 5 6 1 2 3	4 5 6 1 2 3
5	5 6 1 2 3 4	3 4 5 6 1 2
6	6 1 2 3 4 5	2 3 4 5 6 1

【図10】

図 10

ハザード外番号	物理論理変換部	論理物理変換部
	物理ディスク番号	テープ/ハザード番号
	1 2 3 4 5 6	1 2 3 4 5 6
1	1 . 3 4 5 6	1 . 3 4 5 6
2	2 . 4 5 6 1	6 1 . 3 4 5
3	3 . 5 6 1 2	5 6 1 . 3 4
4	4 . 6 1 2 3	4 5 6 1 . 3
5	5 . 1 2 3 4	3 4 5 6 1 .
6	6 . 2 3 4 5	. 3 4 5 6 1

【図11】

図11

ハリテイル外番号	物理論理変換部	論理物理変換部
	物理デイスリ番号	デ-タ/ハリテイル番号
	1 2 3 4 5 6	1 2 3 4 5 6
1	1 . 3 4 5 2	1 6 3 4 5 .
2	2 . 4 5 3 1	6 1 5 3 4 .
3	3 . 5 4 1 2	5 6 1 4 3 .
4	4 . 5 1 2 3	4 5 6 1 3 .
5	5 . 1 2 3 4	3 4 5 6 1 .
6	1 . 2 3 4 5	1 3 4 5 6 .

【図12】

図12

ハリテイル外番号	物理論理変換部	論理物理変換部
	物理デイスリ番号	デ-タ/ハリテイル番号
	1 2 3 4 5 6	1 2 3 4 5 6
1	1 . 3 4 5 2	1 6 3 4 5 .
2	2 . 4 5 3 1	6 1 5 3 4 .
3	3 . 5 4 1 2	5 6 1 4 3 .
4	4 . 5 1 2 3	4 5 6 1 3 .
5	5 . 1 2 3 4	3 4 5 6 1 .
6	1 . 2 3 4 5	1 3 4 5 6 .

【図13】

図13

ハリテイル外番号	物理論理変換部	論理物理変換部
	物理デイスリ番号	デ-タ/ハリテイル番号
	1 2 3 4 5	6 1 2 3 4 5 6
1	1 . 3 . 5 2	1 6 3 . 5 .
2	2 . 4 . 3 1	6 1 5 3 . .
3	3 . 5 . 1 2	5 6 1 . 3 .
4	4 . 5 . 2 3	. 5 6 1 3 .
5	5 . 1 . 3 4	3 . 5 6 1 .
6	1 . 2 . 4 5	1 3 . 5 6 .

【図14】

図14

パリティレイト番号	物理論理変換部	論理物理変換部
	物理ディスク番号	データ/パリティ番号
	1 2 3 4 5 6	1 2 3 4 5 6
1	1 . 3 . 4 2	1 6 3 5 . .
2	2 . 4 . 3 1	6 1 5 3 . .
3	3 . 4 . 1 2	5 6 1 3 . .
4	4 . 1 . 2 3	3 5 6 1 . .
5	2 . 1 . 3 4	3 1 5 6 . .
6	1 . 2 . 4 3	1 3 6 5 . .

【図15】

図15

パリティレイト番号	物理論理変換部	論理物理変換部
	物理ディスク番号	データ/パリティ番号
	1 2 3 4 5 6	1 2 3 4 5 6
1	1 . 3 4 5 2	1 6 3 4 5 .
2	2 . 4 5 3 1	6 1 5 3 4 .
3	3 . 5 4 1 2	5 6 1 4 3 .
4	4 . 6 1 2 3	4 5 6 1 . 3
5	5 . 1 2 3 4	3 4 5 6 1 .
6	1 . 2 3 4 5	1 3 4 5 6 .

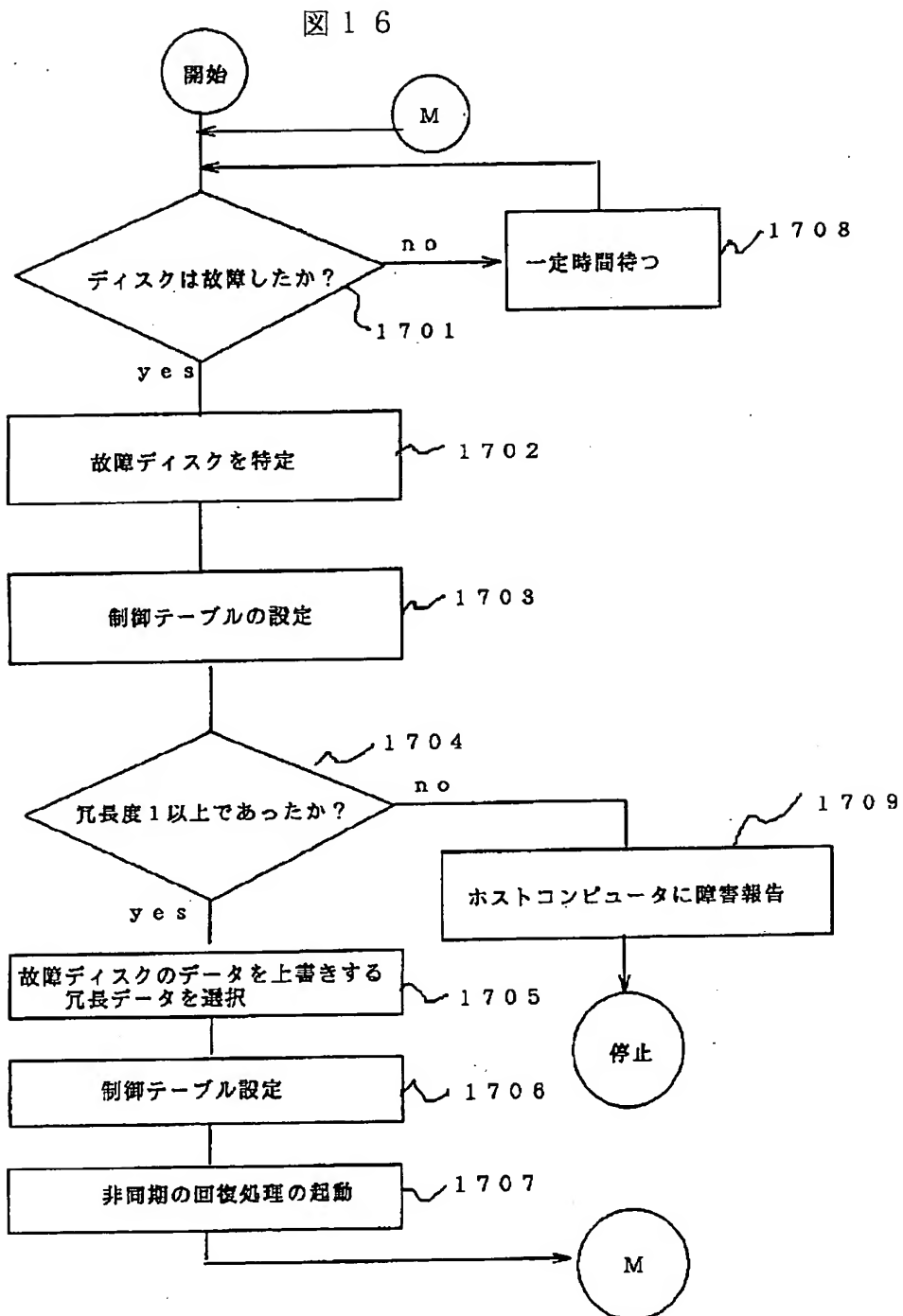
【図27】

図27

	ディスクアクセス回数		
	正常時	故障時	本発明
故障ディスクのデータリード データライト	$\begin{array}{c} \diagup \\ \diagdown \end{array}$	n $n+m-1$	1 $2m$
正常ディスクのデータリード データライト	1 $2(m+1)$	1 $2(m+1)$	1 $2m$

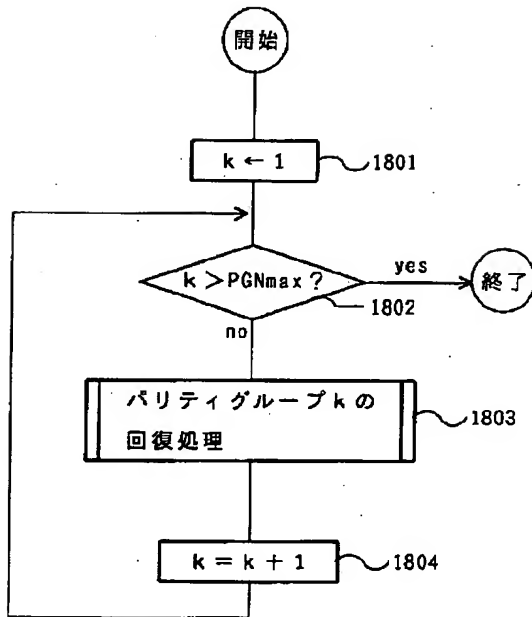
(注) データディスク n 台、パリティディスク m 台の構成を仮定する。
故障はディスク1台の故障を仮定する。

【図16】



【図17】

図 1 7



【図28】

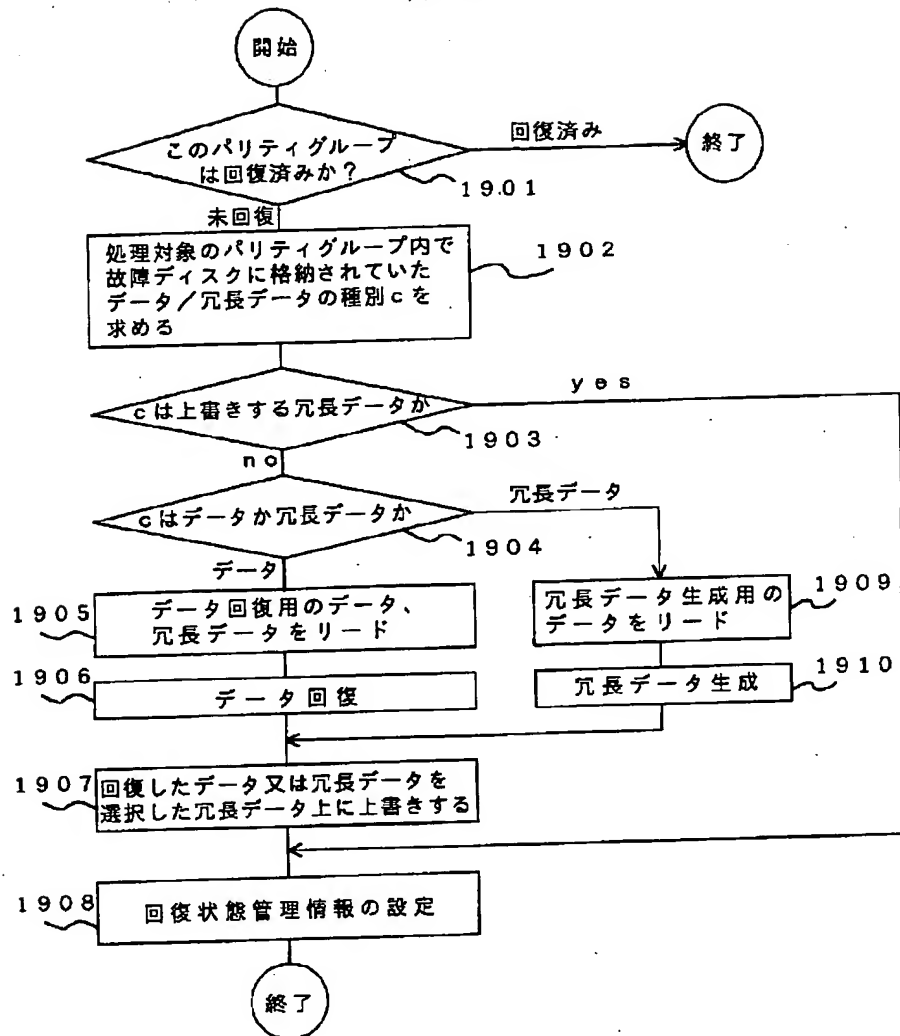
図 2 8

	ディスクアクセス回数		
	正常時	故障時	本発明
故障ディスク内のデータリード	/	n	1
データライト		$n+m-1$	$2m$
対応冗長データが 故障ディスク内にある	/	1	1
データライト		$2m$	$2m$
対応冗長データが 正常ディスク内にある	1	1	1
データライト	$2(m+1)$	$2(m+1)$	$2m$

(注) パリティグループとして、 $(n+m)$ 台のディスクを用いてデータを n 台のディスクに格納し、冗長データを m 台のディスクに格納するものを仮定する。
故障はディスク 1 台の故障を仮定する。

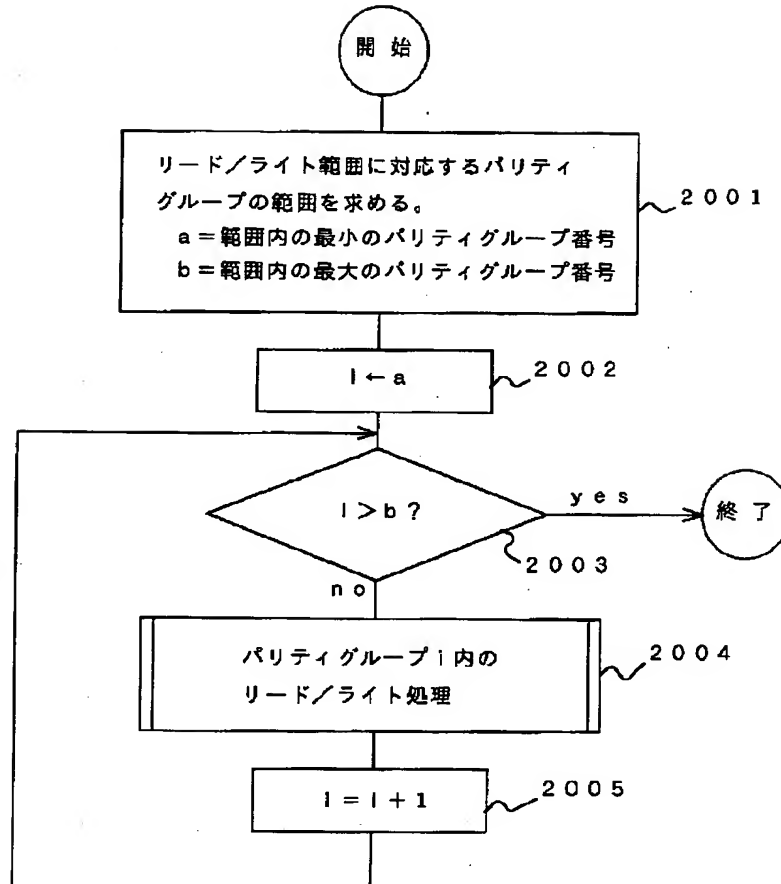
【図18】

図 18



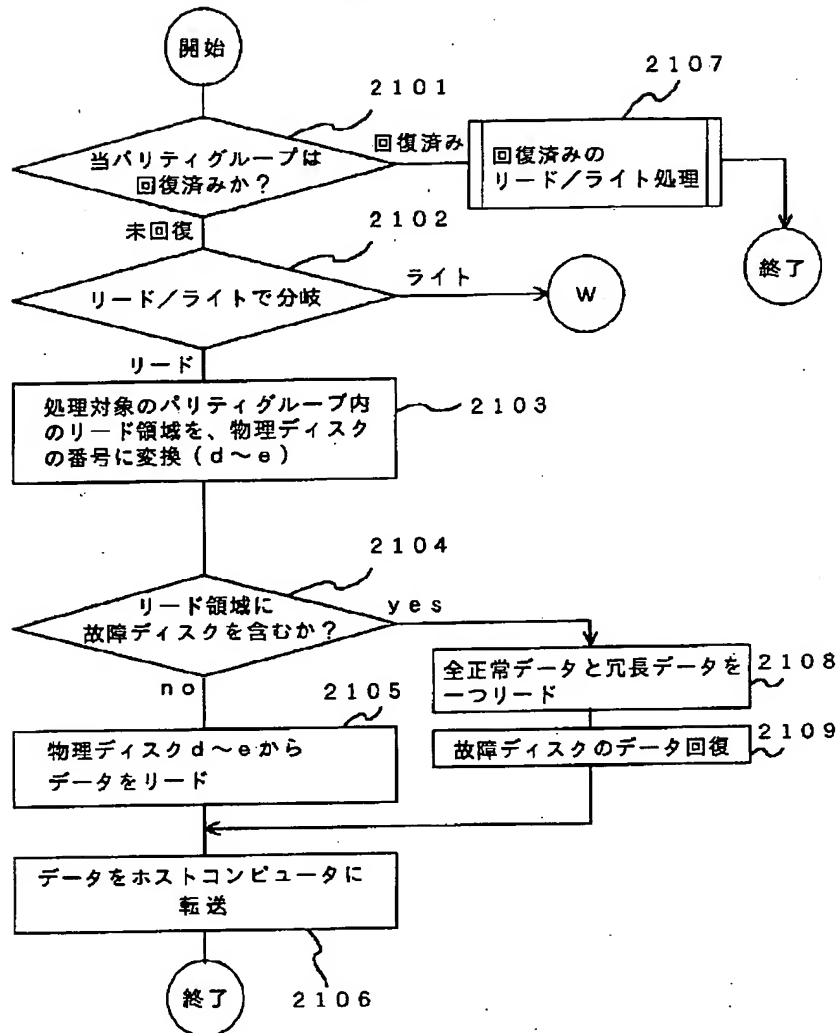
【図19】

図 1 9



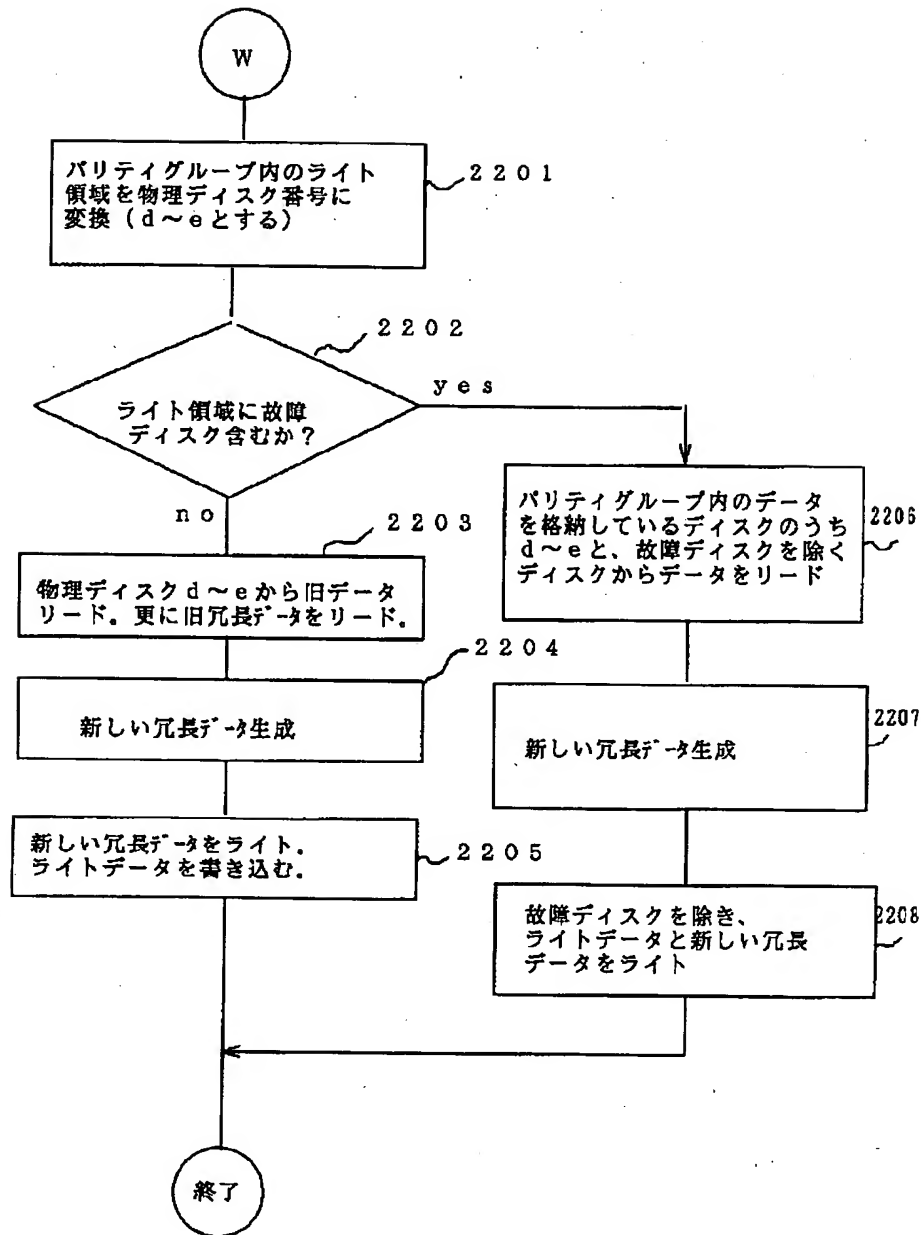
【図20】

図20



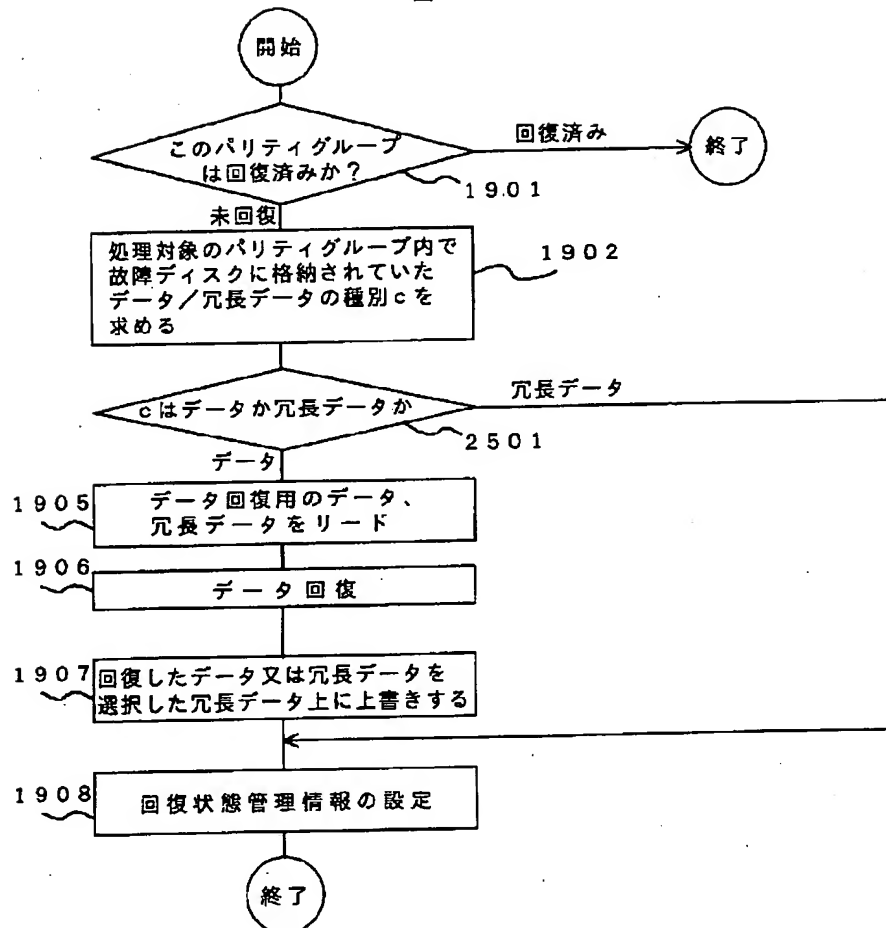
【図21】

図 2 1



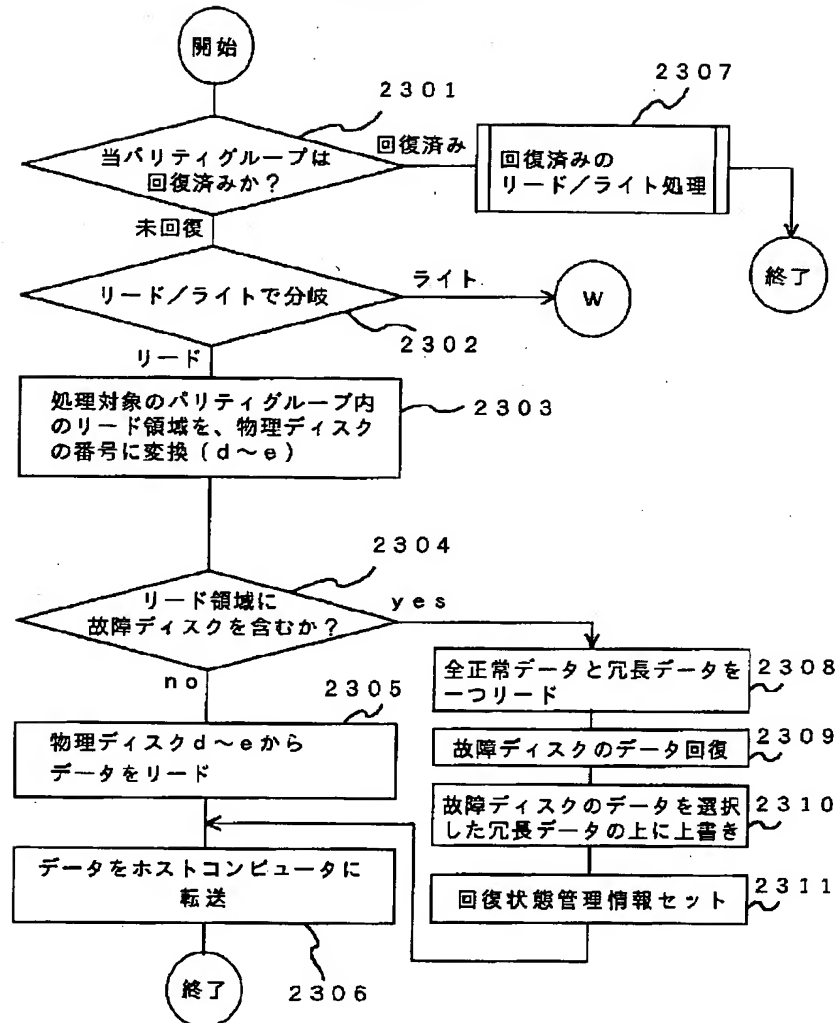
【図22】

図 2 2

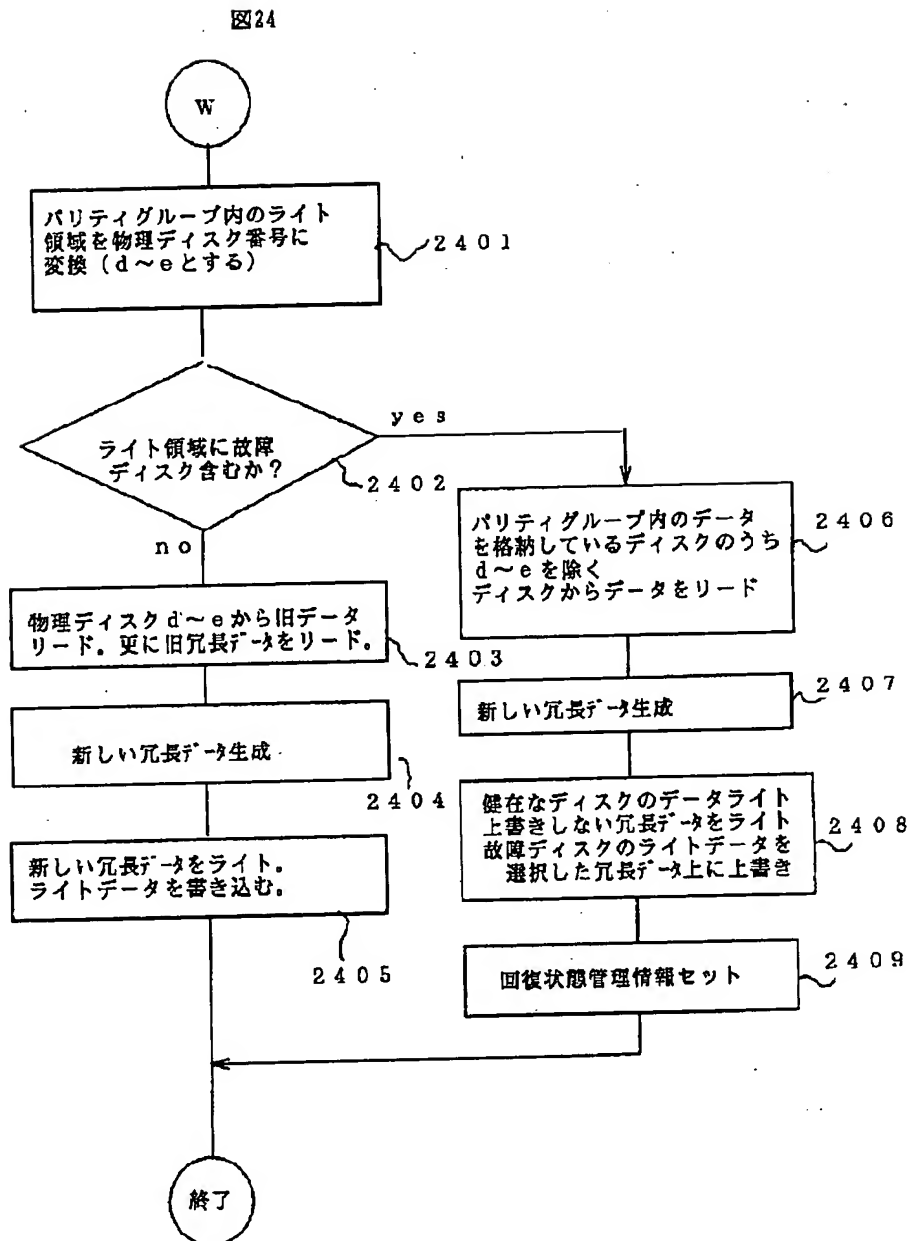


【図23】

図 2 3

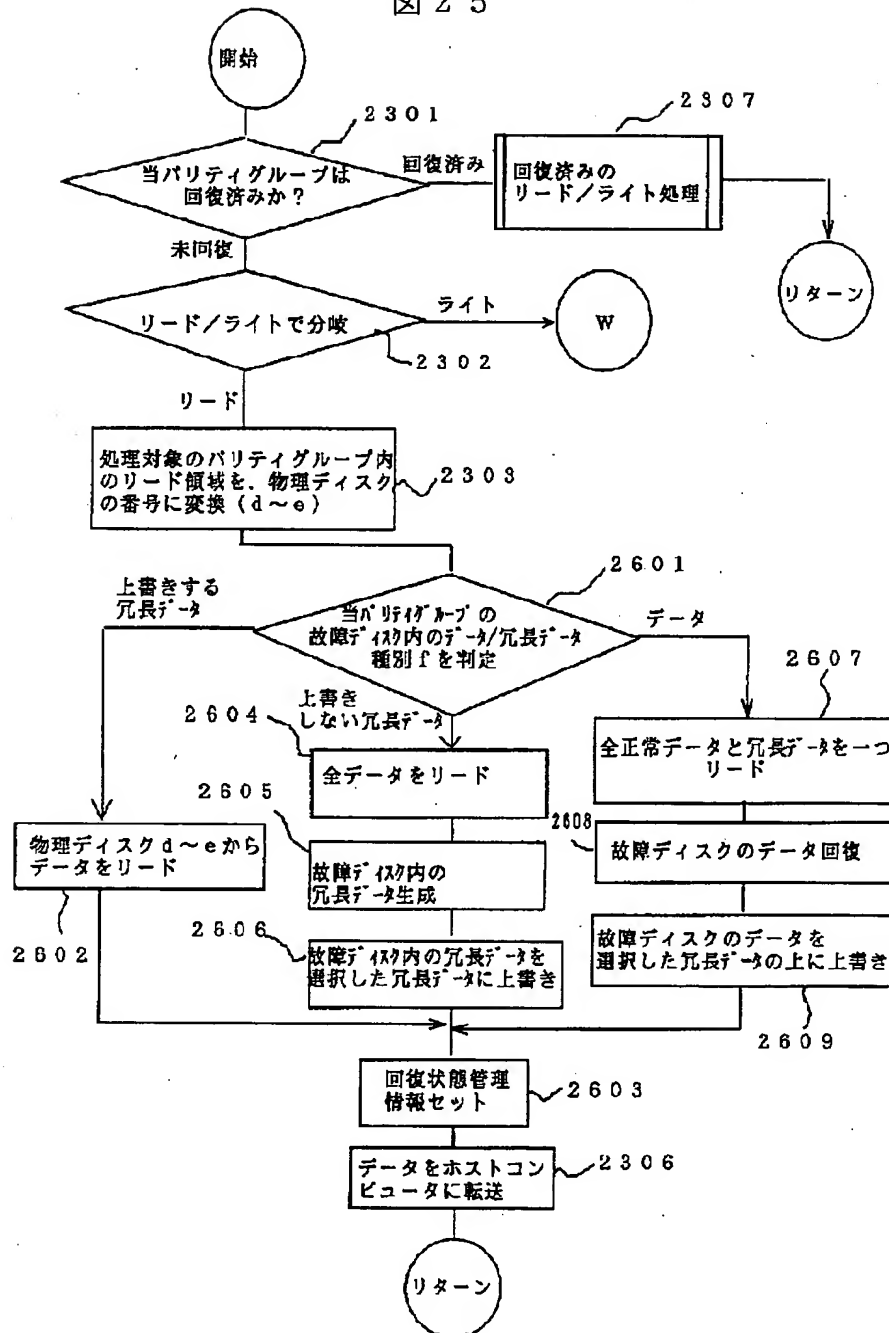


【図24】



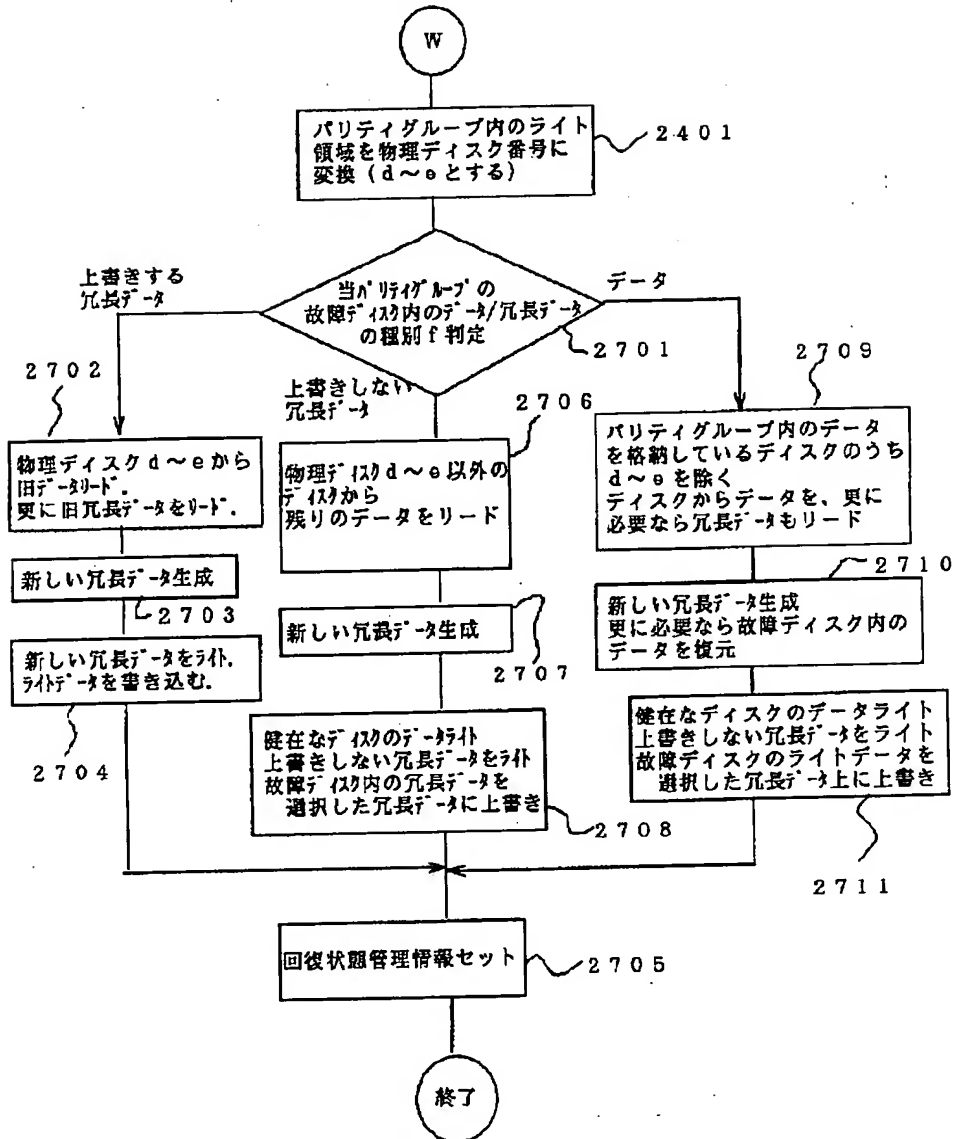
【図25】

図 25



【図26】

図 26



フロントページの続き

(72)発明者 吉田 稔
 神奈川県小田原市国府津2880番地 株式会
 社日立製作所ストレージシステム事業部内

THIS PAGE BLANK (USPTO)